



Midterm Exam

Department: CS

Course Name: Machine Learning

Course Code: CS467 / SCS467

Instructor(s): Dr. Hanaa Bayomi

Name: _____

Date: 22/11/2022

Duration: 1 hour

Total Marks: 20

تعليمات هامة

- حيازة التليفون المحمول ممنوع داخل لجنة الامتحان بغير حالة غش تستوجب العقوب وإذا كان ضروري الدخول بالمحظوظ فلوضع مغلق في العقرب.
- لا يسمح بدخول سماعة الأذن أو الليزروت.
- لا يسمح بدخول أي كتب أو ملازم أو أوراق داخل اللجنة والمختلفة تعتبر حالة غش.

Question 1

[11 marks]

- Answer the following Questions:

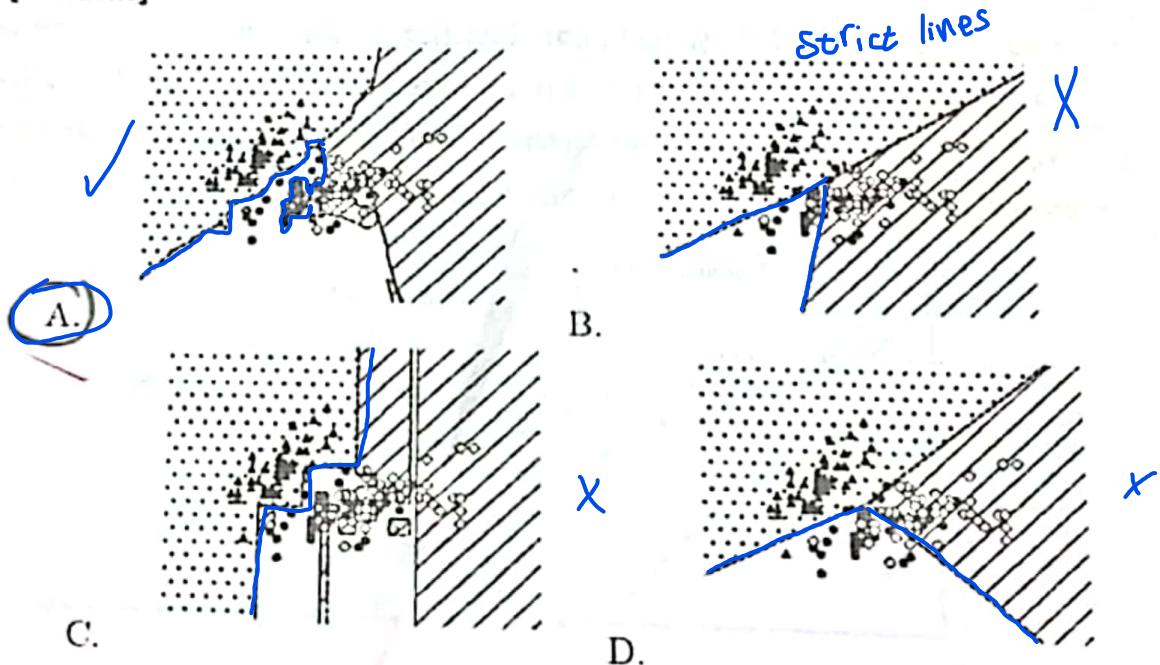
1. We need re-estimate probabilities (smoothing) in Naïve Bayes classifier. [1 mark]

Yes, because features are independent

Yes, to avoid zero conditional probability problem, which happens when no example contains one of the feature values provided in the input (X_i).

2. Which of the following decision boundaries is most likely to be generated by a k-NN?

Why? [2 marks]



The decision boundary generated by k-NN is typically non-linear and depends on the local structure of the data points.

Because it is interested in closed points and works on k-NN

- For the following dataset, calculate the error rate for each KNN classifier. [2 marks]

(b) How many samples are incorrectly classified on the test data for Class 1 (+) and Class 2 (*)? [1 mark]

* 2
+ 1 * 2

[9 marks]

Question 2

You are stranded on a deserted island. Mushrooms of various types grow wildly all over the island, but no other food is anywhere to be found. Some of the mushrooms have been determined as ~~poisonous~~ and others as ~~not~~ (determined by your former companions trial and error). You are the only one remaining on the island. You have the following data to consider.

	IsPoisonous	IsSmelly	IsSalty	IsA Math	IsSpoiled	Class
B	0	0	1	0	0	
C	1	1	0	1	0	
D	1	0	0	1	1	
E	0	1	1	0	1	
F	0	0	1	1	1	
G	0	0	0	1	1	
H	1	1	0	0	1	
U	1	1	1	1	?	
V	0	1	1	1	?	
W	1	1	0	1	?	

You know whether or not mushrooms A through H are poisonous, but you do not know about U through W. For the first couple of questions, consider only mushrooms A through H.

a) What is the entropy of IsPoisonous? [1.5 mark]

$$P_0 = \frac{2}{8} \Rightarrow P_1 = \frac{5}{8}$$

$$\text{Entropy} = -\frac{3}{8} \log\left(\frac{3}{8}\right) - \frac{5}{8} \log\left(\frac{5}{8}\right)$$

$$= 0.954$$

b) Which attributes should you choose as the root of decision tree? [3 mark]

$$H(\text{root}) = 0.954$$

Heavy

$$= 0 \rightarrow 3 \text{ yes}, 2 \text{ no}$$

$$E(H_{v=0}) = -\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5} = 0.97$$

$$= 1 \rightarrow 2 \text{ yes}, 1 \text{ no}$$

$$E(H_{v=1}) = -\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3} = 0.92$$

$$G(\text{root, heavy}) = 0.954 - \frac{5}{8} \times 0.97 - \frac{3}{8} \times 0.92 = 3.2 \times 10^{-3}$$

Smelly

$$\rightarrow 2 \text{ no}, 3 \text{ yes}$$

$$E(Sm=0) = -\frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5} = 0.71$$

$$E(Sm=1) = -\frac{1}{3} \log \frac{1}{3} - \frac{2}{3} \log \frac{2}{3} = 0.918$$

$$Q(S; Sm) = 0.954 - \frac{5}{8} \log \frac{5}{8} - \frac{3}{8} \log \frac{3}{8} = 3.2 \times 10^{-3}$$

Sp

$$\rightarrow 2 \text{ no}, 3 \text{ yes}$$

$$\rightarrow \text{Same as previous}$$

$$E(r, Sp) = 3.2 \times 10^{-3}$$

Smeth

$$\rightarrow 2 \text{ no}, 2 \text{ yes}$$

$$E(Sm) = 1$$

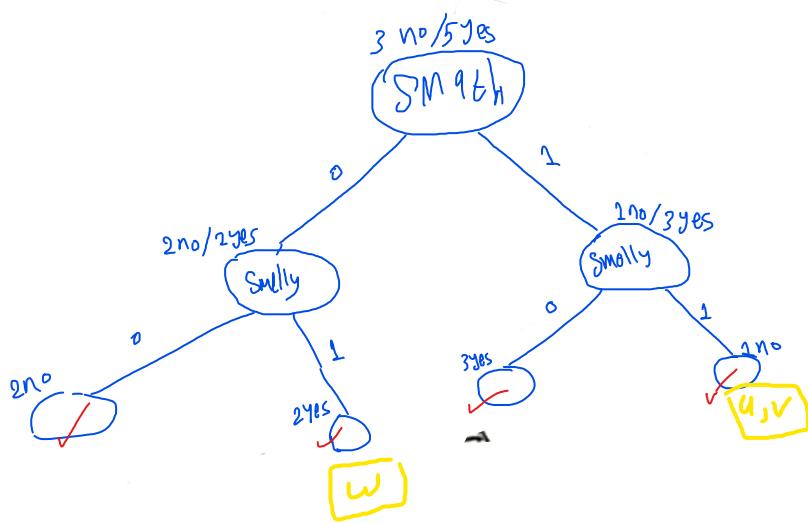
$$\rightarrow 1 \text{ no}, 3 \text{ yes}$$

$$E(r, Sm) = -\frac{1}{4} \log \frac{1}{4} - \frac{3}{4} \log \frac{3}{4} = 0.811$$

$$G(r, Sm) = 0.954 - \frac{5}{8}(1) - \frac{3}{8}(0.811) = 0.0488$$

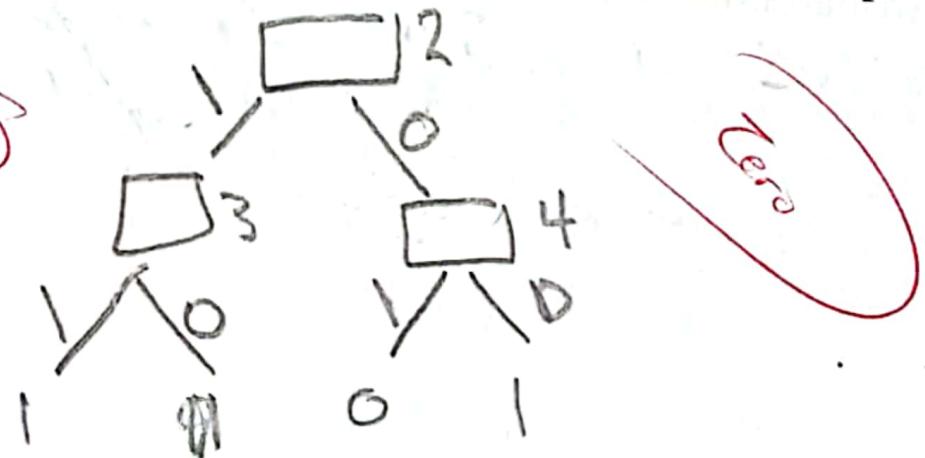
X Highest Gain

c) Build decision tree to classify mushrooms as poisonous or not? [3 marks]

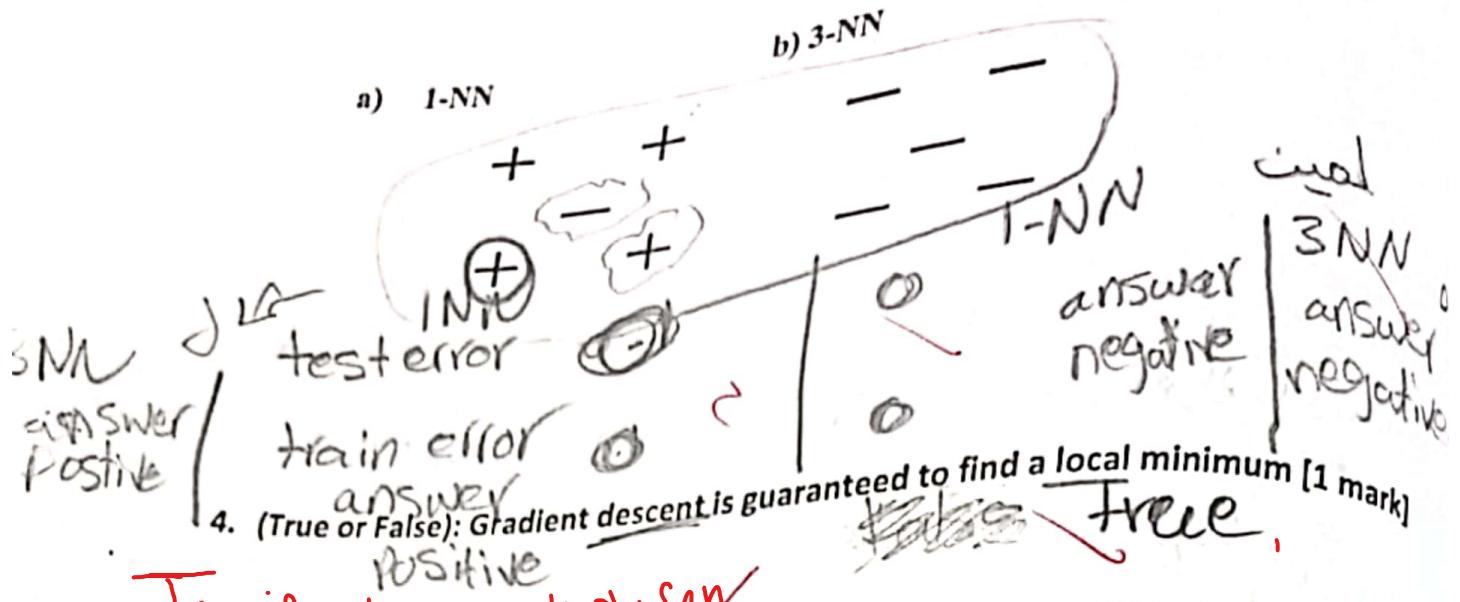


d) Classify mushrooms U,V and W using this decision tree as poisonous or not? [1.5 marks]

$U, V \rightarrow \text{No}$
 $W \rightarrow \text{Yes}$



Good
Dr. Hanaa B

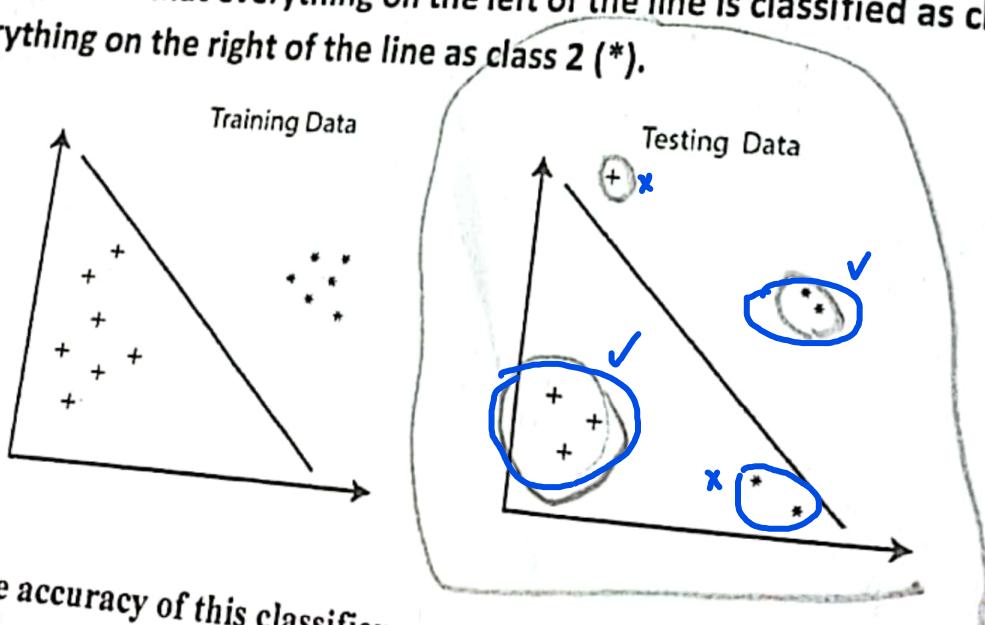


True, if α is correctly chosen

5. Suppose you run gradient descent for linear regression for 100 iterations with a learning rate 0.01. You observe that the training loss (sum of squares error) is increasing after every iteration. How would you explain this observation? And further, what changes would you make to the set up for the gradient descent to converge to a solution? [2 Marks]

~~learn rate is small, need to increase~~
~~data is not normalization, make data noisy~~
Scale by mean

6. You are given the following linear classifier (shown on both training and testing data), with data belonging to class 1 (represented with '+') or class 2 (represented with '*'). Note that everything on the left of the line is classified as class 1 (+), and everything on the right of the line as class 2 (*).



What is the accuracy of this classifier on test data for Class 1 (+) and Class 2 (*)? [2 marks]

$$\frac{T_R + T_H}{J^2} = \frac{\text{Correct}}{\text{total}} = \frac{5}{8}$$

* $\frac{T_R + T_H}{J^2} = \frac{\text{Correct}}{\text{total}} = \frac{5}{8}$

A 3. $\frac{\text{Correct}}{\text{total}} = \frac{5}{8}$